

# Patterns of shared signatures of recent positive selection across human populations

Kelsey Elizabeth Johnson<sup>1</sup> and Benjamin F. Voight<sup>2,3,4\*</sup>

Signatures of recent positive selection often overlap across human populations, but the question of how often these overlaps represent a single ancestral event remains unresolved. If a single selective event spread across many populations, the same sweeping haplotype should appear in each population and the selective pressure could be common across populations and environments. Identifying such shared selective events could identify genomic loci and human traits important in recent history across the globe. In addition, genomic annotations that recently became available could help attach these signatures to a potential gene and molecular phenotype selected across populations. Here, we present a catalogue of selective sweeps in humans, and identify those that overlap and share a sweeping haplotype. We connect these sweep overlaps with potential biological mechanisms at several loci, including potential new sites of adaptive introgression, the glyoxalase locus associated with malarial resistance and the alcohol dehydrogenase cluster associated with alcohol dependency.

Positive selection is the process whereby a genetic variant rapidly increases in frequency in a population due to the fitness advantage of one allele over the other. Recent positive selection has been a driving force in human evolution, and studies of loci targeted by positive selection have uncovered potential adaptive phenotypes in recent human evolutionary history (for example, refs 1–3). One observation that has emerged from scans for positively selected loci is that these signatures often overlap across multiple populations, localized to discrete locations in the genome<sup>4–6</sup>. Sequencing data available from diverse human populations provide an opportunity to characterize the frequency that overlapping signatures share a common, ancestral event—and potentially a common selective pressure. Identifying shared selective events would be of fundamental interest, highlighting loci and traits important in recent history across the globe.

Functional annotations across the human genome can also potentially connect variants targeted by selection with candidate genes and an associated mechanism. For example, the influx of expression quantitative trait loci (eQTL) across many tissue types<sup>7</sup>, and inferred regions of ancient hominin introgression<sup>8</sup>, now provide a richer foundation to investigate the potential biological targets under selection. While identifying the causal variant at a site of positive selection is notoriously difficult, if single nucleotide polymorphisms (SNPs) on a selected haplotype are associated with changes in expression of a nearby gene, this information could help attach the signature to a potential gene and molecular phenotype.

Here, we focus on the detection of genomic signatures compatible with selection on a newly introduced mutation that has not yet reached fixation (that is, a hard, ongoing sweep) to explore their distribution across populations and spanning the genome. We performed a scan for positive selection using the integrated haplotype score (iHS) on 20 populations from four continental groups from phase 3 of the 1000 Genomes Project (1KG)<sup>9</sup>. We found that 88% of sweep events overlapped across multiple populations, correlating

with population relatedness and geographic proximity; 59% of overlaps were shared (that is, a similar sweeping haplotype was present) across populations; and 29% of overlaps were shared across continents. We connect these multi-population sweep overlaps with potential mechanisms at (1) the glyoxalase cluster (GYPF, GYFD and GYFE), where we observe sweeps across all four continental groups in a region associated with malarial resistance; (2) sweeps across African populations at the X chromosome gene *DGKK*, implicated in hypospadias in males; (3) a sweep shared in European populations tagged by a coding variant in the gene *MTA20R*, which is associated with homocysteine levels and a multitude of additional traits; (4) two putative regions of adaptive introgression from Neandertals; and (5) the alcohol dehydrogenase (*ADH*) cluster, where a sweep in Africa is associated with alcohol dependence in African Americans.

## Results and discussion

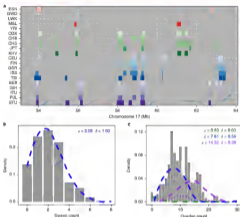
Our broad objective was to apply the iHS to identify and characterize the distribution of selective sweeps across human populations and the genome. Towards this end, we first evaluated the distribution of iHS when applied to whole-genome sequencing data.

**A correction to iHS adjusting for local, low recombination rates.** The iHS was conceptualized for population data ascertained for common genetic variation<sup>10</sup>, and may not be fully calibrated for sequencing data that include rare variation. To examine the score in more detail, we initially applied the iHS to genome sequencing data obtained from 1KG (Methods). We observed an excess of SNPs tagging strong iHS signals at lower derived allele frequencies (Supplementary Fig. 1a) in a frequency range where iHS is not expected to have substantial power<sup>11</sup> (note that a negative iHS indicates extended haplotype homozygosity on the derived relative to the ancestral allele). We observed a negative correlation between the number of populations in an overlap and the local recombination

<sup>1</sup>Genetics and Gene Regulation Program, Cell and Molecular Biology Graduate Group, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA. <sup>2</sup>Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA. <sup>3</sup>Department of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA. <sup>4</sup>Institute for Translational Medicine and Therapeutics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA.

\*e-mail: [voight@perelman.med.upenn.edu](mailto:voight@perelman.med.upenn.edu)





**Fig. 3 | Overlapping sweeps tend to cluster in the genome.** **a**, An example of a 10-Mb window on chromosome 17 with multiple overlaps across many populations. See Supplementary Table 2 for population codes. **b**, The distribution of sweep interval counts in 10-Mb windows across the genome for a single population (GWE). The histogram plots the observed counts, and the blue dashed line is the best-fit Poisson distribution. **c**, The distribution of sweep overlaps across two or more populations in 10-Mb windows across the genome. The histogram plots the observed counts, and the dashed lines represent the results of Poisson mixture modeling. The best-fit model was the three-component model shown here:  $\lambda$ , Poisson distribution rate;  $\hat{\lambda}$ , mixture proportion.

the presence of genes, background selection or local recombination rates. These hotspots could be targeting specific genomic regions, or additional genomic features not assessed here.

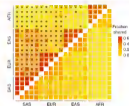
**Complex patterns of sweep sharing across populations and continents.** We next sought to identify selective sweeps that are potentially shared across populations, that is, where the putative sweeping haplotype is similar across populations. Sharing could occur in several ways, including a common ancestral event occurring before population divergence that persisted to the present day, or via gene flow of advantageous alleles between populations. To characterize haplotype similarity across populations at our genomic intervals tagged by unusual dHS, we utilized the program fastPHASE<sup>24</sup>. Using a hidden Markov model, fastPHASE models the observed distribution of haplotypes as mixtures of  $K$  ancestral haplotypes, allowing us to map a sweep tag SNP to an ancestral haplotype jointly across multiple populations at once without arbitrarily choosing a physical span to build a tree of haplotypes or otherwise measure relatedness (Methods).

Overall, out of 1,603 intervals shared across populations, 321 (20%) were shared across continents, frequently between Europe and South Asia, consistent with observed lower genetic differentiation relative to other continental comparisons (Supplementary Table S5). Indeed, consistent with our previous analysis using all intervals,  $F_{ST}$  predicted the fraction of sweep overlaps that were shared between a pair of populations ( $\beta = -2.60$ ,  $\text{s.e.} = 0.34$ ,  $P = 9.1 \times 10^{-16}$ ;

Supplementary Fig. 3). Though more closely related populations have a higher fraction of shared sweeps, they also have more total unshared sweeps. This relationship could be due to shared selective pressures in nearby populations, false negatives in our sharing analysis or a combination of both.

To determine whether the observed extent of sweep sharing was unusual, we applied our fastPHASE haplotype labeling procedure to matched random sites across the genome. For all within-continent population pairs, and all but 4 of 75 Eurasian between-continent pairs, the degree of sweep sharing was higher than the background rate (Fig. 4 and Supplementary Table 6), suggesting that the sweep sharing we observed was not driven purely by haplotype similarities across closely related populations. The number of populations in a shared sweep was inversely correlated with the length of the shared haplotype (Spearman's  $\rho = -0.23$ ,  $P < 1.2 \times 10^{-5}$ , Supplementary Fig. 4), which is compatible with more widely shared sweeping haplotypes being broken down by recombination over time. There was a borderline significant correlation between the nucleotide diversity of a sweeping haplotype and the number of populations sharing a sweep (Spearman's  $\rho = 0.048$ ,  $P = 0.041$ , Supplementary Fig. 5), though this weak correlation was not significant when we looked at diversity between populations within continents.

Though the majority of between-continent shared sweeps were found across non-African populations, we did observe examples of shared sweeps between African and non-African populations. In total, 9.4% of observed sweep overlaps between African and

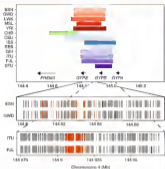


**Fig. 4 | Enrichment of shared sweeps across population pairs.** Squares below the diagonal represent the null fraction of overlaps shared across population pairs, from randomly placed overlaps across the genome. Squares above the diagonal represent the observed fraction of sweep overlaps shared for each population pair. Squares are marked with an asterisk if the observed fraction shared was significantly higher than the null distribution. Populations are arranged alphabetically within continental groups by population code, top to bottom, right to left.

non-African population pairs were called as shared (491 total), compared with 4.0% of control overlaps (99% confidence interval: 3.7–4.4%). For example, on chromosome 1 at ~67 Mb, a sweeping haplotype shared across African and European populations fell in a cluster of cytochrome P450 genes (Supplementary Fig. 6 and Supplementary Note).

**Shared and overlapping sweeps in a region implicated in malaria resistance.** With a catalogue of shared and overlapping selective sweeps in hand, we next aimed to identify specific regions of sweep sharing that connected the interval to a gene, pathway or phenotype when considered alongside additional genomic annotations. With a sweep overlap across thirteen populations from all four continental groups, the glyco-phosphatase gene cluster (*GYPB*, *GYPF*, *GYPH*) came to our attention for its repeated targeting by positive selection and its previous implication in malaria resistance (Fig. 5). This genomic region has been noted as a target of positive selection in humans<sup>22–24</sup>, and as a target of ancient balancing selection shared between humans and chimpanzees<sup>25</sup>. In the IBS and South Asian populations, the sweep appeared to be on a shared haplotype, while the African, CHB and CEU populations each had unique sweeping haplotypes (Fig. 5). This complex locus contains a segmental duplication, making mapping and phasing of short-read data difficult. However, we observed residual unusual dBS in the surrounding region, and in an dBS scan of only those variants passing the 100× ‘strict’ mask. We identified multiple signatures of positive selection on distinct haplotypes in all four continental groups (Supplementary Table 5), with some linked to one or more potentially causal variant(s), either coding or structural (Supplementary Note). The frequency and diversity of apparent adaptive pressures at this locus underscores the role of selection on host–pathogen interactions over recent and longer evolutionary timescales in modern humans, and the potential importance of this locus in particular in that process.

**Intersection of signatures of positive selection with the GWAS catalogue.** Previous work has indicated an enrichment of extreme dBS at genome-wide association studies (GWAS) signals for autoimmune diseases<sup>26</sup>, and we hypothesized that this or other traits might be enriched for GWAS signatures linked to our signatures



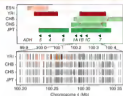
**Fig. 5 | Signatures of positive selection at the *GYP* locus on chromosome 4.** We observed signatures of positive selection in 13 populations at the *GYP* locus, including at least one population from each studied continental group. The top panel shows the sweep intervals of populations with sweeps at this locus, and the positions of the *GYP* genes. The bottom panels show the sweeping haplotypes for two African (CEU, GBR) and two South Asian (ITU, PUA) populations with shared sweeps. The grey tick marks in each population's row indicate the presence of a derived allele on the sweeping haplotype (most common in that population), with a black tick indicating the position of each population's SNP with the most extreme dBS value. Also shown in orange are the significant sQTL for *GYPF* (light orange) or both *GYPB* and *GYPF* (dark orange) in linkage disequilibrium with these populations' shared haplotype (Linkage disequilibrium measure  $r^2 > 1$ ). The sQTL for *GYPB* and *GYPF* are from whole blood.

of positive selection. We tested for enrichment of linked GWAS SNPs overall and for specific traits (autoimmune disorders, height and schizophrenia) by generating random SNP sets from the HapMap3 variant set (Methods), and identifying the number of sweeps in each population linked to a GWAS variant compared with the random SNP sets. We saw no clear, compelling evidence of enrichment for GWAS SNPs overall or any of the traits tested. In total, 186 sweep tag SNPs from all 20 populations (out of 11,653, 1.6%) were in strong linkage disequilibrium with at least one genome-wide significant GWAS SNP ( $r^2 \geq 0.8$ , Supplementary Table 7), compared with a mean of 193 in our random SNP sets (95% confidence interval 168–220). However, this intersection did identify candidates for a potential phenotype under selection at some loci. In one example, a sweep overlap across all five African populations falls at the gene *DGKK* on the X chromosome (Supplementary Note and Supplementary Fig. 7). Variants in this gene have been associated in Europeans with hypopigmentation<sup>27</sup>, a prevalent birth defect of ectopic positioning of the opening of the urethra in males. A second example occurred at the *MTHFR* gene on chromosome 1, where a non-synonymous variant (A222V, rs1801133) has been extensively studied for its association with homocysteine levels<sup>28</sup>. A sweep overlap at this locus with three European populations (CEU, GBR, IBS) and JPT was called as shared across all four populations (Supplementary Note and Supplementary Fig. 8).

**Evidence for adaptive introgression from Neandertals in non-African populations.** Examples of positive selection on introgressed genetic variation have shown that positive selection acted on genetic variation from ancient hominins at some loci. While some of these examples are confined to a single population (for example, *EPAS1* in Tibetans<sup>21</sup>), most are common across multiple populations<sup>22,23</sup>, and thus we hypothesized that a subset of our shared sweeps could be examples of adaptive introgression. We identified 141 candidate sweeps in partial to strong linkage disequilibrium with inferred introgressed haplotypes<sup>24</sup> ( $r^2 \geq 0.6$ , Methods and Supplementary Table 8), including previously described adaptive targets such as the *MTA2* locus in East Asians<sup>25</sup> and *OAS1* in Europeans<sup>26</sup>. We did not observe an overall enrichment of these introgressed haplotypes in our iHS intervals ( $P=0.58$ ,  $\chi^2$  test, Methods), suggesting that introgression alone was not predictive of an unusual iHS signature. Of these 141 loci, we illustrate two candidate sweeps shared across multiple populations (Supplementary Note): (1) a shared sweep between Europeans and South Asians on chromosome 3 near *CTN4*, a non-coding RNA primarily expressed in the testis (Supplementary Fig. 9); and (2) a sweep at  $\sim 41$  Mb on chromosome 1, where all five South Asian populations have evidence of an introgressed haplotype at low to moderate frequency (18–30%) (Supplementary Fig. 10).

**Overlapping and shared sweeps enriched in the ethanol oxidation pathway.** We next sought to explore possible biological pathways targeted by shared selective events. As a large fraction of causal variants under positive selection are potentially non-coding<sup>27–29</sup>, we hypothesized that regulatory variation in the form of eQTL could indicate a potential causal, functional variant and/or gene target. We identified genes with cis eQTL from all tissue types in the GTEx V6p dataset that were linked with shared sweeps ( $r^2 \geq 0.9$ ) and tested for over-representation of biological pathways in this set of genes using ConsensusPathDB<sup>30</sup>. Excluding the human leukocyte antigen (HLA) genes (Methods), the most significant pathway was ethanol oxidation ( $P=2 \times 10^{-4}$ , false discovery rate  $q$  value=0.047), with seven of ten genes included in our shared sweeps gene set. This pathway includes the *ADH* gene cluster, which contains a previously described East Asian selective event targeting rs1229944-T (ref.<sup>31</sup>), a derived non-synonymous variant in *ADH1B* associated with increased *ADH1B* enzyme activity<sup>32</sup> and decreased risk of alcohol dependence in East Asians<sup>33</sup>. A recent report also found evidence for an independent selective event for the variant rs1229944-T in Europeans<sup>34</sup>. Within this region, we observed the East Asian sweep, and independent sweeping haplotypes in the YRI and ESN populations (Fig. 6).

As GWAS have identified genetic variation in the *ADH* locus associated with alcohol dependence<sup>35</sup>, we tested whether these associations were linked to the sweeping haplotype. In the YRI sweep interval spanning *ADH1B*, the derived allele of the leading iHS SNP (rs12639833-T, iHS=5.133) was significantly associated with decreased risk for alcohol dependence in African Americans<sup>35</sup>. This alcohol dependence GWAS in African Americans identified independent associations at a non-synonymous variant in *ADH1B* (rs1266702) and a synonymous variant in *ADH1C* (rs1241494)<sup>35</sup>. The sweep tag in YRI is in perfect linkage disequilibrium with the *ADH1C* lead variant rs1241494 (in LD lies in an intron of *ADH1C* and is a significant eQTL for *ADH1C* (esophagus mucosa) and *ADH4* (esophagus muscularis, skeletal muscle)). Several other SNPs in strong linkage disequilibrium with rs12639833 have extreme negative iHS, are eQTLs for increased *ADH1C* and *ADH4* expression (including in the liver)<sup>36</sup>, and were significantly associated with decreased risk for alcohol dependence in African Americans (Supplementary Table 9). Alcohol dehydrogenase oxidizes ethanol to acetaldehyde, a process that is thought to occur primarily in the liver<sup>37</sup>. These data suggest a similar mechanism is at play in individuals



**Fig. 6 | Signatures of positive selection at the *ADH* locus on chromosome 4.** The top panel shows the sweep intervals of populations with sweeps at this locus, and the positions of the seven *ADH* cluster genes. The bottom panel shows the sweeping haplotypes of the four populations with sweeps in this region, with grey tick marks indicating the derived allele present on the most common sweeping haplotype in that population. The black tick marks indicate the position of the SNP with the most extreme iHS in each population. For the YRI population, the positions of significant *ADH4* and *ADH1C* eQTL in subcutaneous adipose tissue (light orange), GWAS SNPs from ref.<sup>35</sup> (dark orange), and SNPs that are eQTL for both genes and are GWAS SNPs (red) in linkage disequilibrium with YRI's tag SNP ( $r^2 > 0.9$ ) are shown.

of West African ancestry as in East Asians, where the selected allele increases *ADH* enzyme activity<sup>32</sup>, resulting in an adverse physical response to alcohol consumption<sup>38</sup>, and reduced risk for alcohol dependence<sup>33</sup>. Taken collectively, these patterns suggest that (1) alcohol oxidation pathways broadly have been subject to recent positive selection in humans, (2) that genes in this pathway have been repeatedly targeted, with multiple alleles segregating at these sites, (3) the selective pressure appears to operate across the major continental groups included in this study, and (4) sweeping haplotypes at the *ADH* locus tag functional variation associated with protection against alcohol dependence.

In summary, we identified overlapping and shared signatures of positive selection across human populations, using a modified version of the iHS statistic. We observed more extreme iHS in sequencing data compared with SNP array genotype data, which could be a consequence of more rapid decay of homozygosity on unselected haplotypes due to the presence of rare variants. We found that closely related populations are more likely to share sweeping haplotype signatures, though we identified examples of sharing across genetically distant populations. These loci immediately raise questions of how these examples arose, whether by gene flow after divergence or a common ancestral event. Though only a small amount of gene flow between African and non-African populations is thought to have occurred since their divergence, the introduction of an adaptively advantageous allele at very low frequency could lead to the signature we observed. But in considering the collection of putative shared sweeps we highlighted here, it seems apparent that each locus is unique, segregating patterns of genetic variation suggestive of a range of compatible (and potential quite complex) population models that could help explain these data. Future work to infer the potential scenarios leading to shared sweeps will likely require modeling of individual regions to elucidate the evolutionary history of specific events. We also found that the rate of sweep overlaps is not uniform across the genome, but in some locations overlaps cluster together, contributing to the complexity of the underlying sweep signatures in these regions. These features made identifying the tag SNP for a sweep and calling sharing between sweep overlaps difficult in these regions. That said, we hope that our catalogue of

unusually long haplotypes shared across human populations will help to elucidate genes—and ultimately phenotypes—that are still evolving across the wide range of environments human have experienced in recent history

## Methods

**dHS scan.** We downloaded phased genotype files for phase 1 of the 1000 Genomes Project from the 1KG FTP server (<http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110501/>). These data were converted to VCFs (variant call files, and filtered to include only biallelic single nucleotide variants (including indels) with a minor allele frequency (MAF) greater than 1%. A fine-scale recombination map was downloaded from the 1000 FTP server ([http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20110501/cont\\_recombination\\_rates/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20110501/cont_recombination_rates/)), and scaled to units of  $\rho = 4N_e\mu$ , where  $\rho$  is the recombination rate per generation for each population. ELOcore population size ( $N_e$ ) was estimated for each population by calculating nucleotide diversity ( $\pi$ ) as a sliding window (1000 kb) across the genome, and estimating  $N_e$  from the median values of  $\pi$  ( $N_e = \pi/\mu$ , where  $\mu$  is the mutation rate per generation). Ancestral alleles were identified using the human-chimp-mouse alignment from Ensembl (sourced from [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/biallelic/analysis\\_results/supporting/ancestral\\_alignments/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/biallelic/analysis_results/supporting/ancestral_alignments/)). SNPs were filtered for only those where the ancestral allele was supported by both the chimp and mouse alignments.

Standardized dHS were calculated using WHAMM ([doi:10.1038/nature11044](http://doi.org/10.1038/nature11044)), using a modified version of dHS calculation code (v1.1) that increased speed of the calculation, and initially standardized by derived allele frequency as described in the original dHS paper<sup>1</sup>, with 50 allele frequency bins. In the final standardization, we binned ancestral SNPs into 500 bins (50 allele frequency bins  $\times$  10 local recombination rate bins), or 100 bins for chromosomes X (50 allele frequency bins  $\times$  2 local recombination rate bins). These standardization files are available in Supplementary Table 1.

Regions of the genome putatively undergoing recent hard sweeps—what we refer to as the main test and dHS intervals—were identified by counting the number of dHSs with  $|dHS| > 2$  in 100 kb windows (windows exceeding by one SNP that is overlapping windows). We took the mean of the top 1% of windows, by the total number and by fraction of dHSs with  $|dHS| > 2$  in the window, as our intervals. We performed this interval calling separately for each of the 10 populations included in this study. The SNP we used to label (that is, tag) each sweep interval was identified as the SNP with the most extreme dHS, and the sweep frequency as the tag SNP-derived allele frequency if the dHS was less than zero, and ancestral allele frequency if dHS was greater than zero. We limited our analysis of individual sweep calls to those with a tag SNP of MAF  $> 15\%$ , to focus on signatures unlikely to have extreme dHS due to very low frequency.

**Neutral simulations.** We performed neutral simulations of a population with a CIBI population-like demographic history using the forward simulation software SLiM<sup>27</sup> (v 2.4.1). We simulated chromosomes with the following demographic model<sup>27</sup>—an ancestral population size of 13,000 individuals with a mutation rate of  $1 \times 10^{-8}$  was fixed at for 130,000 generations, then a bottleneck reduced the population size to 2,000 individuals at 11,000 generations, followed by exponential growth at a rate of  $1.002502$ , and finally 116 diploids were sampled at generation 14,500. We simulated 85,116 regions at a recombination rate of  $r_c = 1 \times 10^{-8}$  or  $1 \times 10^{-9}$ , and calculated standardized dHS for all variants. Only variants in the central 100 kb were used for the comparison in Supplementary Fig. 2, to maximize the number of high-frequency ( $>1\%$  MAF) variants that reached the end of their haplotype within the simulated chromosomes (for a total of 159,487 variants at  $r_c = 1 \times 10^{-8}$  and 116,718 variants at  $r_c = 1 \times 10^{-9}$ ). The mean of standardized dHS for variants in an allele frequency bin was compared between the difference recombination rates with the Mann-Whitney test (wilcox.test), and the variances were compared with the F test (var.test) in R<sup>28</sup>.

**Sweep overlaps.** To identify sweep overlaps, we compared the dHS intervals for each population and identified regions of the genome where two or more populations had a sweep overlap. We calculated the fraction of sweep overlaps for each population pair as the mean of the fraction of sweep intervals in one population that overlap with a sweep interval in the second population (that is,  $\text{fraction in population A} \times \text{fraction in population B} \div 2$ ). We estimated  $P_{ij}$  for each pair of populations across all variants ( $n = 2,617,340$ ) in the 1000 Genomes VCF files on chromosome 2 using the Weir and Cockerham estimator implemented in VCFtools<sup>29</sup> (v 0.1.16)<sup>30</sup>. We performed linear regression on  $\text{fraction of overlap} = P_{ij}$  for each population pair, and estimated the standard error of the slope using a block bootstrap for unequal group size<sup>31</sup>, deleting one population pair group (for example, EUR versus EUR, or AFR versus EUR) for each comparison.

**Runs across the genome.** To assess the rate of sweep intervals across the genome, we subdivided the genome into 50 kb non-overlapping windows ( $n = 200$ ) in total and counted the number of sweep intervals for each individual population, and the number of overlaps across two or more populations, in each window. To assess the sweep intervals called for each population were independent,

we merged adjacent sweep intervals into one interval if three tag SNPs were in modest linkage disequilibrium or greater ( $r^2 > 0.4$ ). If a sweep interval or overlap spanned two windows, we counted it once in the window with more than half of its physical distance. We used all 1000 Genomes phase 1 gene annotations ([http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/working/20110501/cont\\_recombination\\_rates/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/working/20110501/cont_recombination_rates/)) and haplotype selection estimates (S statistics) were downloaded from <http://www.1000genomes.org/analysis/selection/> ([doi:10.1038/nature11044](http://doi.org/10.1038/nature11044)). We performed forward regression with the medians of 8 and local recombination rates in our windows as potential predictors of overlap count in each window, along with the gene count in the window. The final model included only median local  $R^2 = 0.048$ ,  $P = 5.2 \times 10^{-5}$ , as median  $R^2$  explained no additional variance. For Poisson count modeling of the overlap rate, we fit mixtures of independent Poisson distributions to the data by maximizing the negative log likelihood with the nonlinear mixture function (nlm) in R<sup>28</sup>. We compared mixture models by calculating the Bayesian information criterion and performing a likelihood ratio test.

**Identifying sweeping haplotypes with fastPHASE.** For each sweep overlap, we identified the physical region spanning all tag SNPs, and an additional 1 kb to either side. We ran fastPHASE on this region, using the  $\pi$  option to identify each HGS population as a subpopulation,  $R$  to indicate known haplotypes and  $\text{Pop}$  to output cluster probabilities for each individual at each SNP. We tested a range of values of  $K$  (number of haplotype clusters) and  $T$  (number of random expectation maximization algorithm starts) on a subset of sweep overlaps, and found broadly similar results across the range (Supplementary Fig. 11). We used  $K = 10$  clusters and  $T = 10$  for all overlaps in the final analysis. From the output cluster probabilities, we identified the sequence of haplotype clusters for each SNP position in each individual as the most likely haplotype cluster at each SNP. We then identified the haplotype cluster sequence of all chromosomes carrying the selected tag allele, and the most common of these to be the reference sweeping haplotype sequence.

To identify a pair of populations as 'shared', we required an identical reference haplotype sequence to span the selected tag allele in both populations. To form shared clusters, we grouped together all populations that were called as shared with at least one other population. To calculate the null rate of haplotype sharing across population pairs, we selected random regions of the genome of the same size (within 10 kb), distant to the current gene (within 50 kb), and local recombination rate (within an order of magnitude of  $\rho = 4N_e\mu$ ) as our observed sweep overlap regions. For each sweep overlap, we identified 10 simulated windows, for a total of 10,000 regions across the genome (sampling from 15 to 1,500 random overlaps per population pair). We identified tag SNPs for each population in the random regions matching the distance from the other population's tag SNPs and derived allele frequency (within 1%) of the observed overlap. We then ran fastPHASE on the randomly selected regions and performed the shared haplotype calling procedure as for observed-sweep windows described above. To compare the observed fraction of overlaps called as shared to the null haplotype sharing for each pair of populations, we performed 1,000 bootstraps by sampling with replacement the number of observed overlaps from the null. Population pairs where the shared sweep fraction of observed overlaps was higher than the shared fraction of random overlaps for all 1,000 samples are marked with an asterisk in Fig. 4. We performed linear regression on  $\text{fraction of shared overlaps} = P_{ij}$  for each population pair, and estimated the standard error of the slope using a block bootstrap for unequal group size<sup>31</sup>, deleting one population pair group (for example EUR versus EUR, or AFR versus EUR) for each comparison.

We measured the length of the shared sweeping haplotype in each population as the maximum length of the reference shared haplotype, identified as above, in that population. We calculated nucleotide diversity in each population as the average number of pairwise differences between chromosomes carrying the tag allele in the region of the shared haplotype, divided by the length of the shared haplotype in that population. We tested for correlations between the number of populations sharing a sweep and the mean shared haplotype length and mean nucleotide diversity using Spearman's rho in R (cor.test). We also tested, in each population individually, the correlation between nucleotide diversity of the shared haplotype in that population and the number of populations in the shared sweep.

**qTL linked to sweep haplotypes.** To connect shared sweeps to potential causal genes, we utilized the GTEx v6p qTL dataset<sup>32</sup> downloaded from the GTEx portal (<http://www.gtexportal.org/>). For each population with dHSs, we identified linkage disequilibrium priors ( $r^2 > 0.8$ ) calculated in the same population) within 1 kb of the sweep interval, and intersected these SNPs with all significant GTEx qTLs from all tissue types. qTLs in the GTEx V6p dataset were identified among a subset of mostly white individuals (84.1%), with a smaller fraction of African Americans (13.7%). For sweep overlaps that were called as shared, we identified a shared SNP set as the intersections of linkage disequilibrium priors sets for all populations in a shared group. We created a gene hit of all genes with qTLs from any tissue that intersected with shared SNP sets, excluding HLA genes. We chose to exclude HLA genes, owing to its genomic complexity and its enrichment for signatures of recent positive selection. To test for enrichment of this gene set with biological pathways,

we used gene representations analysis of all pathway databases in *ConsensusPathDB* (<http://cpath.molgen.mpg.de/>) with the background set of all genes.

**Intersection of sweeps with Nonneutral haplotypes, Nonneutral P/WAS and GWAS SNPs.** We downloaded the Nonneutral haplotypes call reported in ref. 11 from [http://haploref.jax.washington.edu/veroot\\_et\\_al\\_2016\\_release\\_data/](http://haploref.jax.washington.edu/veroot_et_al_2016_release_data/). This dataset contains inferred introgressed Nonneutral haplotypes on the autosomes of the non-African individuals from 100 phase 3. To calculate linkage disequilibrium between introgressed haplotypes and sweep tag SNPs, we pooled overlapping haplotypes across individuals and created a genotype of 0 or 1 based on presence or absence of the overlapping introgressed haplotype in each individual. We then calculated linkage disequilibrium between this genotype or absence genotype and the tag SNPs within 1 kb of the introgressed haplotype separately for each population, then considered haplotypes with  $r^2 > 0.5$  with sweep tag SNPs as candidates for adaptive introgression. To examine potential enrichment of introgressed haplotypes in linkage disequilibrium with sweep tag SNPs, we compared the fraction of introgressed haplotypes in linkage disequilibrium with sweep tag SNPs to the fraction of distance and frequency matched SNPs in linkage disequilibrium with sweep tag SNPs ( $r^2 > 0.5$ ) using a  $\chi^2$  test. We downloaded the Nonneutral phenotype-wide association studies (P/WAS) data at [https://phenoscan.org/veroot/et\\_al/](https://phenoscan.org/veroot/et_al/) and intersected all reported associations with variants in strong linkage disequilibrium ( $r^2 \geq 0.9$ ) with each sweep tag SNP in each population.

We downloaded the GWAS catalogue from <https://www.ebi.ac.uk/gwas/> 22 October 2016. We identified all genome-wide significant associations ( $P < 5 \times 10^{-8}$ ) in strong linkage disequilibrium ( $r^2 \geq 0.9$ ) with each sweep tag SNP in each population. To test for enrichment of GWAS variants generally and of specific phenotype classes, we performed permutation tests with 10,000 random SNPs from the HapMap3 variant set from <http://cpath.molgen.mpg.de/hapmap3/>, 30 matched for allele frequency and distance to gene with the GWAS variants of interest. We then compared the empirical distribution of intersection of these matched SNPs sets with the sweep tag SNPs and protein to the number of observed GWAS intersections. We control for potentially linked GWAS variants, we simply counted the number of sweeps in each population that intersected a GWAS or control set variant.

**Indels, structural variants and annotations.** Indels were not included in our HSI scan, but could be the causal variant on a sweeping haplotype. To identify candidates for causal indels, we calculated linkage disequilibrium with sweeping SNPs for all indels in the 1000 Genomes phase 3 VCF files within 1 kb of the sweep interval in each population. To identify potential functional coding variants among indels and SNPs on sweeping haplotypes, we used ANNOVAR to annotate coding variation<sup>33</sup>. In the glycoprotein regions, we tested for linkage disequilibrium between sweep tag SNPs and structural variant calls from phase 3 of the 1000 ([http://1000genomes.ebi.ac.uk/docs/phase3/merged\\_vc\\_mphs/24/February2010\\_release/](http://1000genomes.ebi.ac.uk/docs/phase3/merged_vc_mphs/24/February2010_release/)).

**Life Science Reporting Summary.** Further information on experimental design is available in the Life Science Reporting Summary.

**Code availability.** WDRAMD, WDR4 and HSI calculator v1.5 are available at <http://consortium.stat.su.se/wdramd/hsiscan.html>. Custom code used to process and analyse output from the HSI scan and fastPANS is available at <https://github.com/hsiscan/hsiscan>.

**Data availability.** The standardized HSI values for all 26 populations in 100 phase 3 are available at [http://consortium.stat.su.se/data/hsiscan/A\\_HSIvalues.txt.gz](http://consortium.stat.su.se/data/hsiscan/A_HSIvalues.txt.gz).

Received: 28 June 2017; Accepted: 11 January 2018;

Published online: 19 February 2018

## References

- Tufekci, S. A. et al. Convergent adaptation of human lactase persistence in Africa and Europe. *Nat. Genet.* **38**, 36–40 (2006).
- Kanfer, T. G. et al. Modeling recent human evolution in mice by sequences of a selected H2B variant. *Cell* **152**, 699–705 (2013).
- Pomajbí, M. et al. Greenlandic Inuit show genetic signatures of diet and climate adaptation. *Science* **349**, 1345–1347 (2015).
- Pollard, J. D. et al. Signals of recent positive selection in a worldwide sample of human populations. *Genome Res.* **19**, 824–837 (2009).
- Cong, Q. et al. The role of geography in human adaptation. *PLoS Genet.* **3**, e1000200 (2007).
- Moripala, M. et al. Shared and unique components of human population structure and genome-wide signals of positive selection in South Asia. *Am. J. Hum. Genet.* **88**, 731–744 (2011).
- Liu, X. et al. Detecting and characterizing genomic signatures of positive selection in global populations. *Am. J. Hum. Genet.* **93**, 861–881 (2013).
- Agar, P. et al. Genetic effects on gene expression across human tissues. *Nature* **526**, 206–213 (2017).
- Veroot, S. & Akay, I. M. Resurrecting surviving nonneutral lineages from modern human genomes. *Science* **345**, 1017–1021 (2014).
- Azou, A. et al. A global reference for human genetic variation. *Nature* **526**, 48–74 (2015).
- Vogel, B. P., Kondrinski, S., Wu, X. & Pritchard, J. E. A map of recent positive selection in the human genome. *PLoS Biol.* **4**, e72 (2006).
- Salari, P. C. et al. Genome-wide detection and characterization of positive selection in human populations. *Nature* **449**, 913–918 (2007).
- McVicker, G., Gordon, D., Davis, C. & Green, P. Widespread genomic signatures of natural selection in hominid evolution. *PLoS Genet.* **5**, e1000571 (2009).
- Schum, P. & Stephens, M. A fast and flexible statistical model for large-scale population genetic data: applications to inferring missing genotypes and haplotype phase. *Am. J. Hum. Genet.* **78**, 825–844 (2006).
- Joost, L., Wient, R. H. & Corns, D. J. Natural selection on the cytoplasmic surface. *Mol. Biol. Evol.* **30**, 225–239 (2013).
- Wang, H. Y., Tang, H., Shen, C. K. & Wu, C. I. Rapidly evolving genes in humans. I. The glycoproteins and their possible role in creating modern humans. *Mol. Biol. Evol.* **20**, 1795–1804 (2003).
- Ko, W. Y. et al. Effects of natural selection and gene conversion on the evolution of human glycoproteins coding for MNS blood polymorphisms in malaria-endemic African populations. *Am. J. Hum. Genet.* **88**, 741–754 (2011).
- Leffler, E. M. et al. Resistance to malaria through structural variation of red blood cell invasion receptors. *Science* **336**, 1531–1537 (2011).
- Leffler, E. M. et al. Multiple instances of convergent balancing selection shared between humans and chimpanzees. *Science* **318**, 1378–1382 (2007).
- Taj, T. et al. Common risk alleles for inflammatory diseases are targets of recent positive selection. *Am. J. Hum. Genet.* **92**, 523–529 (2013).
- van der Zanden, L. P. M. et al. Common variants in CDK6 are strongly associated with risk of hypogonadism. *Nat. Genet.* **43**, 48–50 (2011).
- Griffin, P. et al. Genome-wide association analyses identify variants in developmental genes associated with hypogonadism. *Nat. Genet.* **46**, 557–563 (2014).
- Park, G. et al. Novel associations of CYP11, MUTE, NR2E3, and CYP11 with plasma homocysteine in a healthy population: a genome-wide evaluation of 13 074 participants in the women's genome health study. *Circ. Cardiovasc. Genet.* **2**, 140–150 (2009).
- van Meulen, J. E. et al. Common genetic loci influencing plasma homocysteine concentrations and their effect on risk of coronary artery disease. *Am. J. Clin. Nutr.* **98**, 688–678 (2013).
- Huerta-Sánchez, E. et al. Altitude adaptation in Tibetans caused by introgression of Denisovan-like DNA. *Nature* **521**, 194–197 (2014).
- Raciono, F., Sanjurjo-Roman, S., Nielsen, R. & Huerta-Sánchez, E. Evidence for adaptive introgression in humans. *Nat. Rev. Genet.* **16**, 159–171 (2015).
- Raciono, F., Manríquez, D. & Huerta-Sánchez, E. Signatures of archaic adaptive introgression in present-day human populations. *Mol. Biol. Evol.* **34**, 296–317 (2017).
- Ding, Q., Hu, Y., Xu, S., Wang, J. & Jin, L. Nonneutral introgression at chromosome 1p31.1 was under positive natural selection in east Asians. *Mol. Biol. Evol.* **34**, 681–695 (2017).
- Jain, A. et al. Adaptively introgressed Neandertal haplotypes at the OAS1 locus functionally impacts innate immune responses in humans. *Genome Biol.* **17**, 244 (2016).
- Grossman, S. R. et al. Identifying recent adaptations in large-scale genomic data. *Cell* **152**, 705–713 (2013).
- Enard, D., Messer, P. W. & Petro, D. A. Genome-wide signals of positive selection in human evolution. *Genome Res.* **24**, 989–995 (2014).
- Yu, F. et al. Population genomic analysis of 962 whole genomes identifies of human reveals natural selection in non-coding regions. *PLoS Genet.* **10**, e1003444 (2014).
- Randev, A., Wierling, C., Lohrke, H. & Wierling, B. ConsensusPathDB—a database for integrating human functional interaction networks. *Nucleic Acids Res.* **35**, D423–D428 (2009).
- Hua, Y. et al. Evidence of positive selection on a class I ADH locus. *Am. J. Hum. Genet.* **60**, 61–65 (2007).
- Fehring, H. J. The genetics of alcohol metabolism: role of alcohol dehydrogenase and aldehyde dehydrogenase variants. *Alcohol Res. Health* **30**, 3–13 (2007).
- Li, D., Zhao, H. & Gelernter, J. Strong associations of the alcohol dehydrogenase 1B gene (ADH1B) with alcohol dependence and alcohol-related medical diseases. *Biol. Psychiatry* **74**, 596–601 (2013).
- Gallay, K. J. et al. Four principal-component analysis reveals convergent evolution of ADH1B in Europe and East Asia. *Am. J. Hum. Genet.* **96**, 436–451 (2015).

38. Gelernter, J. et al. Genome-wide association study of alcohol dependence: significant findings in African- and European-Americans including novel risk loci. *Mol. Psychiatry* **19**, 41–49 (2014).
39. Paulsen, O. E. et al. Large, diverse population cohorts of hEPSCs and derived hepatocyte-like cells reveal functional genetic variation at blood lipid-associated loci. *Cell Stem Cell* **26**, 558–570.e18 (2021).
40. Lee, S. L., Chien, G. Y., Yeh, C. T., Wu, C. W. & Yin, S. J. Functional assessment of human alcohol dehydrogenase family in ethanol metabolism: significance of first pass metabolism. *Alcohol Clin. Exp. Res.* **36**, 1810–1842 (2006).
41. Matsuo, K. et al. *Alcohol dehydrogenase 1 class1A9 polymorphism influences drinking habit independently of aldehyde dehydrogenase 3 class1A7 polymorphism: analysis of 1,299 Japanese subjects*. *Cancer Epidemiol. Biomark. Prev.* **15**, 1809–1815 (2006).
42. Hellen, B. C. & Messer, P. W. Slah 3. Recomb. interactive between genetic simulations. *Mol. Biol. Evol.* **34**, 230–240 (2017).
43. Terhorst, J., Kamm, J. A. & Song, T. S. Robust and scalable inference of population history from hundreds of unphased whole genomes. *Nat. Genet.* **49**, 885–899 (2016).
44. R Development Core Team R: A Language and Environment for Statistical Computing (R Foundation for Statistical Computing, Vienna, 2014).
45. Ihmshagen, P. et al. The mutant cell format and VCPtools. *Bioinformatics* **27**, 1038–1038 (2011).
46. Huang, F. M., T. A., Meyer, E. & Jordan van der, B. Delphi: in pipelines for transcript ex. *Bioinformatics* **9**, 5–8 (1999).
47. Quacken, A. B. & Hall, L. M. RCTtools: a flexible suite of utilities for comparing genomic datasets. *Bioinformatics* **26**, 841–843 (2010).
48. Vernot, B. et al. Increasing Neanderthal and Denisovan DNA from the genomes of Melanesian individuals. *Science* **353**, 201–208 (2016).

49. Simons, C. H. et al. The phenotypic legacy of admixture between modern humans and Neanderthals. *Science* **381**, 732–741 (2018).
50. Wang, K., Lu, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high throughput sequencing data. *Nucleic Acids Res.* **36**, e164 (2008).

**Acknowledgements**

We thank L. Brown and B. Taffes for helpful comments that improved the quality of the manuscript. This work was supported by grants from the National Institutes of Health (MH090304MH1479 to R.P.V., T32CA000026 to K.E.J.) and a fellowship from the Alfred P. Sloan Foundation (202012-087 to R.P.V.).

**Author contributions**

K.E.J. and R.P.V. planned the study. K.E.J. assembled input data and performed the experiments. K.E.J. and R.P.V. interpreted the data and wrote the paper.

**Competing interests**

The authors declare no competing financial interests.

**Additional information**

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41574-018-0479-4>.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).  
Correspondence and requests for materials should be addressed to R.P.V.  
**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.